

1 Generating counterfactual explanations of
2 tumor spatial proteomes to discover
3 therapeutic strategies for enhancing immune
4 infiltration

5 Zitong Jerry Wang^{1*}, Alexander M. Xu², Aman Bhargava¹
6 and Matt W. Thomson^{1*}

7 ¹Division of Biology and Biological Engineering, California
8 Institute of Technology, 1200 E California Blvd, Pasadena, 91125,
9 CA, USA.

10 ²Samuel Oschin Comprehensive Cancer Institute, Cedars-Sinai
11 Medical Center, 8700 Beverly Blvd, Los Angeles, 90048, CA, USA.

12 *Corresponding author(s). E-mail(s): zwang2@caltech.edu;
13 mthomson@caltech.edu;

14 **Abstract**

15 Immunotherapies can halt or slow down cancer progression by activat-
16 ing either endogenous or engineered T cells to detect and kill cancer
17 cells. For immunotherapies to be effective, T cells must be able to infil-
18 trate the tumor microenvironment. However, many solid tumors resist
19 T-cell infiltration, challenging the efficacy of current therapies. Here,
20 we introduce Morpheus, an integrated deep learning framework that
21 takes large scale spatial omics profiles of patient tumors, and combi-
22 nes a formulation of T-cell infiltration prediction as a self-supervised
23 machine learning problem with a counterfactual optimization strat-
24 egy to generate minimal tumor perturbations predicted to boost T-cell
25 infiltration. We applied our framework to 368 metastatic melanoma
26 and colorectal cancer (with liver metastases) samples assayed using
27 40-plex imaging mass cytometry, discovering cohort-dependent, combi-
28 natorial perturbations, involving CXCL9, CXCL10, CCL22 and CCL18
29 for melanoma and CXCR4, PD-1, PD-L1 and CYR61 for colorec-
30 tal cancer, predicted to support T-cell infiltration across large patient

31 cohorts. Our work presents a paradigm for counterfactual-based pre-
32 diction and design of cancer therapeutics using spatial omics data.

33 Introduction

34 The immune composition of the tumor microenvironment (TME) plays
35 a crucial role in determining patient prognosis and response to cancer
36 immunotherapies [1–3]. Immunotherapies that alter the immune composition
37 using transplanted or engineered immune cells (chimeric antigen receptor T
38 cell therapy) or remove immunosuppressive signaling (checkpoint inhibitors)
39 have shown exciting results in relapsed and refractory tumors in hematolog-
40 ical cancers and some solid tumors. However, effective therapeutic strategies
41 for most solid tumors remain limited [4–6]. The TME is a complex mixture of
42 immune cells, including T cells, B cells, natural killer cells, and macrophages,
43 as well as stromal cells and tumor cells [1]. The interactions between these
44 cells can either promote or suppress tumor growth and progression, and ulti-
45 mately impact patient outcomes. For example, high levels of tumor-infiltrating
46 lymphocytes (TILs) in the TME are associated with improved prognosis and
47 response to immunotherapy across multiple cancer types [7, 8]. Conversely,
48 an immunosuppressive TME characterized by low levels of TILs is associated
49 with poor prognosis and reduced response to immunotherapy [9]. Durable,
50 long-term clinical response of T-cell-based immunotherapies are often con-
51 strained by a lack of T-cell infiltration into the tumor, as seen in classically
52 “cold” tumors such as triple-negative breast cancer or pancreatic cancer, which
53 have seen little benefit from immunotherapy [10–12]. The precise cellular and
54 molecular factors that limit T-cell infiltration into tumors is an open question.

55 Spatial omics technologies capture the spatial organization of cells and
56 molecular signals in intact human tumors with unprecedented molecular detail,
57 revealing the relationship between localization of different cell types and tens
58 to thousands of molecular signals [13]. T-cell infiltration is modulated by
59 a rich array of signals within the tumor microenvironment (TME) such as
60 chemokines, adhesion molecules, tumor antigens, immune checkpoints, and
61 their cognate receptors [14]. Recent advances in *in situ* molecular profiling
62 techniques, including spatial transcriptomic [15, 16] and proteomic [17, 18]
63 methods, simultaneously capture the spatial relationship of tens to thousands
64 of molecular signals and T cell localization in intact human tumors with
65 micron-scale resolution. Imaging mass cytometry (IMC) is one such technol-
66 ogy that uses metal-labeled antibodies to enable simultaneous detection of up
67 to 40 antigens and transcripts in intact tissue [17].

68 Recent work on computational methods as applied to multiplexed tumor
69 images have primarily focused on predicting patient-level phenotypes such as
70 survival, by identifying spatial motifs from tumor microenvironments [19–22].

71 These methods have generated valuable insights into how the complex compo-
72 sition of TMEs influences patient prognosis and treatment response, but they
73 fall short of generating concrete, testable hypotheses for therapeutic interven-
74 tions that may improve patient outcomes. Given the prognostic significance of
75 T-cell infiltration into tumors, we need computational tools that can predict
76 immune cell localization from environmental signals and systematically gener-
77 ate specific, feasible tumor perturbations that are predicted to alter the TME
78 to improve patient outcomes.

79 Counterfactual explanations (CFEs) can provide important insight in
80 image analysis applications [23], but have not been applied to multiplexed
81 imaging data. Traditionally, CFEs help clarify machine learning model deci-
82 sions by exploring hypothetical scenarios, showing how the model's interpre-
83 tation would change if a feature in an image were altered slightly [24]. For
84 instance, slight pixel intensity variations or minor edge alterations in a tumor's
85 appearance on an X-ray might lead a diagnostic model to classify the scan
86 differently. Numerous CFE algorithms exist to elucidate a model's decision
87 boundaries and shed light on its sensitivity to specific image features [25]. In
88 multiplexed tissue images where each pixel captures detailed molecular infor-
89 mation, variations in pixel intensity directly correspond to specific molecular
90 interventions. Thus, spatial omics data enables the extension of CFEs from
91 understanding to predicting actionable interventions.

92 In this work, we introduce Morpheus, an integrated deep learning frame-
93 work that first leverages large scale spatial omics profiles of patient tumors to
94 formulate T-cell infiltration prediction as a self-supervised machine learning
95 (ML) problem, and combines this prediction task with counterfactual opti-
96 mization to propose tumor perturbations that are predicted to boost T-cell
97 infiltration. Specifically, we train a convolutional neural network to predict T-
98 cell infiltration using spatial maps of the TME provided by IMC. We then apply
99 a gradient-based counterfactual generation strategy to the infiltration neural
100 network to compute changes to the signaling molecule levels that increase pre-
101 dicted T-cell abundance. We apply Morpheus to melanoma [26] and colorectal
102 cancer (CRC) with liver metastases [27] to discover tumor perturbations that
103 are predicted to support T cell infiltration in tens to hundreds of patients.
104 We provide further validation of ML-based T-cell infiltration prediction using
105 an additional breast cancer data set [28]. For patients with melanoma, Mor-
106 pheus predicts combinatorial perturbation to the CXCL9, CXCL10, CCL22
107 and CCL18 levels can convert immune-excluded tumors to immune-inflamed
108 in a cohort of 69 patients. For CRC liver metastasis, Morpheus discovered
109 two cohort-dependent therapeutic strategies consisting of blocking different
110 subsets of CXCR4, PD-1, PD-L1 and CYR61 that are predicted to improve
111 T-cell infiltration in a cohort of 30 patients. Our work provides a paradigm
112 for counterfactual-based prediction and design of cancer therapeutics based on
113 classification of immune system activity in spatial omics data.

Results

Counterfactual optimization for therapeutic prediction

The general logic of Morpheus (Figure 1A) is to first train, in a self-supervised manner, a classifier to predict the presence of CD8+ T cells from multiplexed tissue images (Figure 1B). Then we compute counterfactual instances of the data by performing gradient descent on the input image, allowing us to discover perturbations to the tumor image that increases the classifier's predicted likelihood of CD8+ T cells being present (Figure 1C). The altered image represents a perturbation of the TME predicted to improve T-cell infiltration. We mask CD8+ T cells from all images to prevent the classifier from simply memorizing T-cell expression patterns, guiding it instead to learn environmental features indicative of T-cell presence.

We leverage IMC profiles of human tumors to train a classifier to predict the spatial distribution of CD8+ T cell in a self-supervised manner. Consider a set of images $\{I^{(i)}\}$, obtained by dividing IMC profiles of tumor sections into local patches of tissue signaling environments, where $I^{(i)} \in \mathbb{R}^{l \times w \times c}$ is an array with l and w denoting the pixel length and width of the image and c denoting the number of molecular channels in the images (Figure 1B). Each image shows the level of c proteins across all cells within a small patch of tissue. From patch $I^{(i)}$, we obtain a binary label $s^{(i)}$ indicating the presence and absence of CD8+ T cells in the patch and a masked copy $x^{(i)}$ with all signals originating from CD8+ T cells removed (see Methods). The task for the model f is to classify whether T cells are present ($s^{(i)} = 1$) or absent ($s^{(i)} = 0$) in image $I^{(i)}$ using only its masked copy $x^{(i)}$. Specifically, $f(x^{(i)}) \in [0, 1]$ is the predicted probability of T cells, and then we apply a classification threshold p to convert this probability to a predicted label $\hat{s}^{(i)} \in \{0, 1\}$. Since we obtain the image label $s^{(i)}$ from the image $I^{(i)}$ itself by unsupervised clustering of individual cells, our overall task is inherently self-supervised.

Given a set of image patches, we train a model f to minimize the following T cell prediction loss, also known as the binary cross entropy (BCE) loss,

$$L = -\frac{1}{N} \sum_{i=1}^N \left[s^{(i)} \log(\hat{s}^{(i)}) + (1 - s^{(i)}) \log(1 - \hat{s}^{(i)}) \right], \quad (1)$$

where

$$\hat{s}^{(i)} = \begin{cases} 1 & \text{if } f(x^{(i)}) \geq p \\ 0 & \text{if } f(x^{(i)}) < p \end{cases} \quad (2)$$

and p is the classification threshold. We select p by minimizing the following root mean squared error (RMSE) on a separate set of tissue sections Ω ,

$$\text{RMSE}^2 = \frac{1}{|\Omega|} \sum_{j \in \Omega} \left| \frac{1}{N_j} \sum_{i=1}^{N_j} s^{(i)} - \frac{1}{N_j} \sum_{i=1}^{N_j} \hat{s}^{(i)} \right|^2. \quad (3)$$

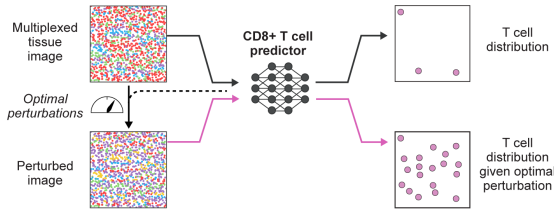
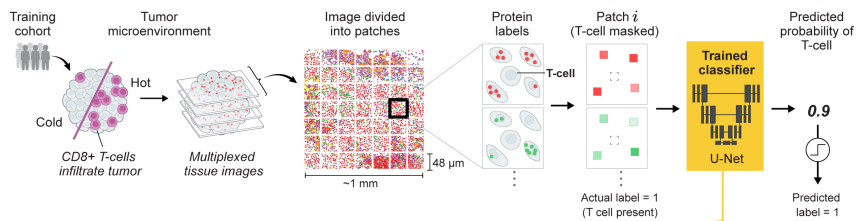
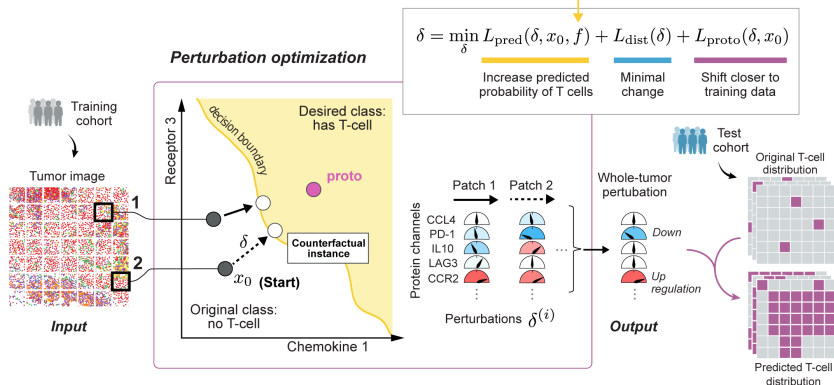
A Overview of Morpheus: a counterfactual optimization framework**B Self-supervised training of T cell localization classifier****C Counterfactual optimization of tissue perturbation**

Fig. 1: An integrated counterfactual optimization framework for discovering therapeutic strategies predicted to drive CD8+ T cell infiltration in human tumors. (A) Overview of the Morpheus framework, which consists of first (B) training a neural network classifier to predict the presence of CD8+ T cells from multiplexed tissue images where CD8+ T cells are masked. (C) The trained classifier is then used to compute an optimal perturbation vector $\delta^{(i)}$ per patch by jointly minimizing three loss terms (L_{pred} , L_{dist} , L_{proto}). The perturbation $\delta^{(i)}$ represents a strategy for altering the level of a small number of signaling molecules in patch $x_0^{(i)}$ in a way that increases the probability of T cell presence as predicted by the classifier. The optimization also favors perturbations that shift the image patch to be more similar to its nearest T-cell patches in the training data, shown as proto. Each perturbation corresponds to adjusting the relative intensity of each imaging channel. Taking the median across all perturbations produces a whole-tumor perturbation strategy, which we assess by perturbing *in silico* tumor images from a test patient cohort and examining the predicted T cell distribution after perturbation.

144 The RMSE is a measure of the differences between the observed and pre-
 145 dicted proportions of T cell patches in a tissue section averaged across a set of
 146 tissues Ω , which we take to be the validation set.

147 We evaluated the performance of various classifiers, including both tradi-
 148 tional convolutional neural networks (CNNs) and vision transformers. In all
 149 cases, we observed similar performance (Table S3). We settled on a U-Net
 150 architecture because of ease of extension of the model to multichannel data
 151 sets. Our U-Net classifier consists of a standard U-Net architecture [29] and
 152 a fully-connected layer with softmax activation (Methods). To increase the
 153 number of samples available for training, we take advantage of the spatial het-
 154 erogeneity of TMEs and divide each tissue image into $48 \mu\text{m} \times 48 \mu\text{m}$ patches
 155 upon which the classifier is trained to predict T cell presence (Methods).

Using our trained classifier and IMC images of tumors, we employ a counter-
 factual optimization method to predict tumor perturbations that enhance
 CD8+ T cell infiltration (Figure 1C). For each image patch $x_0^{(i)}$ lacking CD8+
 T cell, our optimization algorithm searches for a perturbation $\delta^{(i)}$ such that
 our classifier f predicts the perturbed patch $x_p^{(i)} = x_0^{(i)} + \delta^{(i)}$ as having T cells,
 hence $x_p^{(i)}$ is referred to as a counterfactual instance. Ideally, we also want our
 perturbation to be minimal in that it only requires targeting a small number
 of molecule, and realistic in that the counterfactual instance is not far from
 image patches in our training data so we can be more confident of the model's
 prediction. We can obtain a perturbation $\delta^{(i)}$ with these desired properties by
 solving the following optimization problem adopted from [30],

$$\delta^{(i)} = \min_{\delta} L_{\text{pred}}(x_0^{(i)}, \delta) + L_{\text{dist}}(\delta) + L_{\text{proto}}(x_0^{(i)}, \delta), \quad (4)$$

such that

$$\begin{aligned} L_{\text{pred}}(x_0^{(i)}, \delta) &= c \max(-f(x_0^{(i)} + \delta), -p), \\ L_{\text{dist}}(\delta) &= \beta \|\delta\|_1 + \|\delta\|_2^2, \\ L_{\text{proto}}(x_0^{(i)}, \delta) &= \theta \|x_0^{(i)} + \delta - \text{proto}^{(i)}\|_2^2 \end{aligned} \quad (5)$$

156 where $\delta^{(i)}$ is a 3D tensor that describes perturbation made to each pixel of the
 157 patch.

158 The three loss terms in Equation (4) each correspond to a desirable prop-
 159 erty of the perturbation we aim to discover. The term L_{pred} encourages validity,
 160 in that the perturbation increases the classifier's predicted probability of T
 161 cells to be larger than p , so the network will predict the perturbed tissue patch
 162 as having T cells when it previously did not contain T cells. Next, the term
 163 L_{dist} encourages sparsity, in that the perturbation does not require making
 164 many changes to the TME, by minimizing the distance between the original
 165 patch $x_0^{(i)}$ and the perturbed patch $x_p^{(i)} = x_0^{(i)} + \delta$ using elastic net regulariza-
 166 tion. Lastly, the term $\text{proto}^{(i)}$ in the expression for L_{proto} refers to the nearest
 167 neighbour of $x_0^{(i)}$ among all patches in the training set that are classified as

168 having T cells (see [Methods](#)). Thus the term L_{proto} explicitly guides the per-
 169 turbed image $x_p^{(i)}$ to lie close to the data manifold defined by our training set,
 170 making perturbed patches appear similar to what has been observed in TMEs
 171 infiltrated by T cells.

Since drug treatments cannot act at the spatial resolution of individual micron-scale pixels, we constrain our search space to only perturbations that affect all cells in the image uniformly. Specifically, we only search for perturbations that change the level of any molecule by the same relative amount across all cells in an image. We incorporate this constraint by defining $\delta^{(i)}$ in the following way,

$$\delta^{(i)} = \gamma^{(i)} \odot_3 x_0^{(i)}, \quad (6)$$

172 where $\gamma^{(i)} \in \mathbb{R}^c$ defines a single factor for each channel in the image and the
 173 circled dot operator represent channel-wise multiplication, so that within each
 174 channel, the scaling factor is constant across the spatial dimensions of the
 175 image. In practice, we directly optimize for $\gamma^{(i)}$, where $\gamma_j^{(i)}$ can be interpreted
 176 as the relative change to the mean intensity of the j -th channel. However,
 177 given our classifier does have fine spatial resolution, we can search for targeted
 178 therapies such as perturbing only a specific cell type or restricting the per-
 179 turbation to specific tissue locations by changing Equation (6) to match these
 180 different types of perturbation.

181 Taken together, our algorithm obtains an altered image predicted to contain
 182 T cells from an original image which lacks T cells, by minimally perturbing the
 183 original image in the direction of the nearest training patch containing T cells
 184 until the classifier predicts the perturbed image to contain T cells. Since our
 185 strategy may find different perturbations for different tumor patches, we reduce
 186 the set of patch-wise perturbations $\{\delta^{(i)}\}_i$ to a whole-tumor perturbation by
 187 taking the median across the entire set.

188 Convolutional neural networks predict T-cell distribution

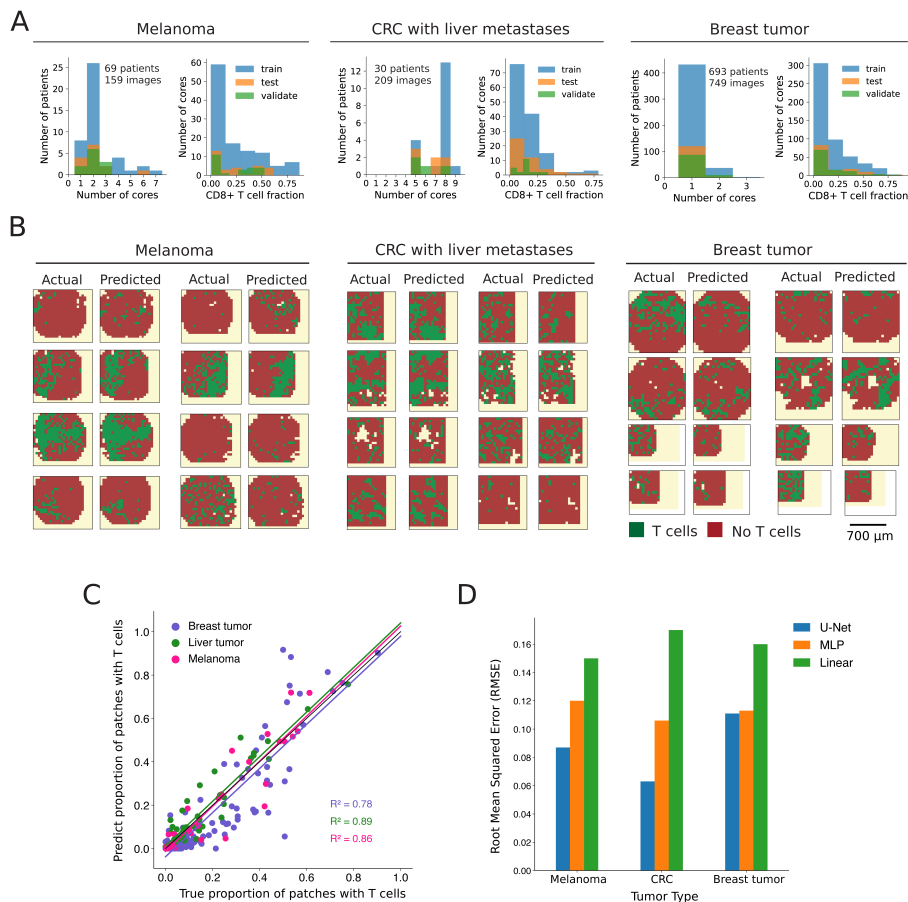


Fig. 2: U-Net classifiers accurately predict T cell distribution in IMC images of melanoma, metastatic liver, and breast tumor. (A) Histograms showing the distribution of tumor cores per patient and CD8+ T cell fractions per core across all three data sets and data splits. (B) Predicted and actual T cell distribution of tissue sections from test cohorts in melanoma, liver tumor, and breast tumor data set. (C) Predicted and true proportion of patches with T cells within a tissue section, each dot corresponds to a tissue section, diagonal black line indicates perfect prediction. (D) The RMSE (Equation (3)) across all (test) tissue sections for three different classes of models.

189 We applied Morpheus to two publicly available IMC data sets of tumors from
 190 patients with metastatic melanoma [26] and colorectal cancer (CRC) with
 191 liver metastases [27] (Figure 2A). We validate the infiltration prediction on an

192 additional breast cancer data set [28]. While this breast cancer data focuses
193 on cell type markers over functional modulators of T-cell infiltration, making
194 it unsuitable for therapeutic prediction, it serves to further validate our ML-
195 based prediction of T-cell infiltration.

196 The melanoma data set [26] was obtained by IMC imaging of 159 tumor
197 cores from 69 patients with stage III or IV metastatic melanoma. Each tis-
198 sue was imaged across 39 molecular channels, consisting of markers for tumor,
199 immune, and stromal cells, as well as 11 different chemokines (RNA) ([Meth-](#)
200 [ods](#)). The CRC data set [27] consists of 209 tissue sections taken from 30
201 patients imaged across 42 channels, including 60 sections from primary CRC
202 tumors, 89 sections CRC metastases to the liver and 60 “healthy” liver sections
203 obtained away from the metastases ([Methods](#)). The breast cancer data set [28]
204 was obtained by IMC imaging of 749 breast tumor cores from 693 patients.
205 The tissues were imaged across 37 channels, consisting of markers for tumor,
206 lymphoid, myeloid and stromal cells ([Methods](#)).

207 For each of the three tumor data sets, we trained a separate U-Net clas-
208 sifier that effectively predicts CD8+ T cell infiltration level in unseen tumor
209 sections ([Methods](#)). The two classifiers trained on melanoma and CRC data
210 sets achieved the best performance with an AUROC of 0.77 and 0.8 respec-
211 tively, whereas the classifier trained on breast tumors achieved a AUROC of
212 0.71 ([Table S2](#)). [Figure 2B](#) shows examples of actual and predicted T cell
213 distributions in tumor sections. For each tissue section of a cancer type, the pre-
214 dictions were obtained by applying the corresponding U-Net classifier to each
215 image patch independently. By visual inspection, our classifiers consistently
216 captures the general distribution of T cells. Comparing the true proportion of
217 T-cell patches in a tissue section against our predicted proportion also shows
218 strong agreement ([Figure 2C](#)). The true proportion of patches with T cells
219 is calculated by dividing the number of patches within a tissue section that
220 contain CD8+ T cells by the total number of patches within that section.
221 We quantify the performance of our U-Nets on the entire test data set using
222 the RMSE (Equation (3)), which represents the mean difference between our
223 predicted proportion and the true proportion per tumor section ([Figure 2D](#)).
224 Our classifiers performs well on liver tumor and melanoma, achieving a RMSE
225 of only 6% and 8% respectively and a relatively lower performance of 11%
226 on breast tumor. Taken together, these results suggest that our classifier can
227 accurately predict the T cell infiltration status of multiple tumor types.

228 In order to gain insight into the relative importance of non-linearity and
229 spatial information in the performance of the U-Net on the T cell classification
230 task, we compared the U-nets’ performance to a logistic regression model (LR)
231 and a multi-layer perceptron (MLP). Both the LR and MLP model are given
232 only mean channel intensities as input, so neither have explicit spatial infor-
233 mation. Furthermore, the LR model is a linear model with a threshold whereas
234 the MLP is a non-linear model. [Figure 2D](#) shows that across all three cancer
235 data sets, the MLP classifier consistently outperforms the logistic regression
236 model, reducing RMSE by 20 – 40% to suggest that there are significant non-
237 linear interactions between different molecular features in terms of their effect

238 on T cell localization. The importance of spatial features on the T cell pre-
239 diction task, however, is less consistent across cancer types. [Figure 2D](#) shows
240 that for predicting T cells in breast tumor, the U-Net model offers negligible
241 boost in performance relative to the MLP model ($< 2\%$ RMSE reduction),
242 whereas for liver tumor, the U-Net model achieved a RMSE 50% lower com-
243 pared to the MLP model. This result suggests that the spatial organization
244 of signals may have a stronger influence on CD8+ T cell localization in liver
245 tumor compared to breast tumor.

246 Applying Morpheus to metastatic melanoma samples

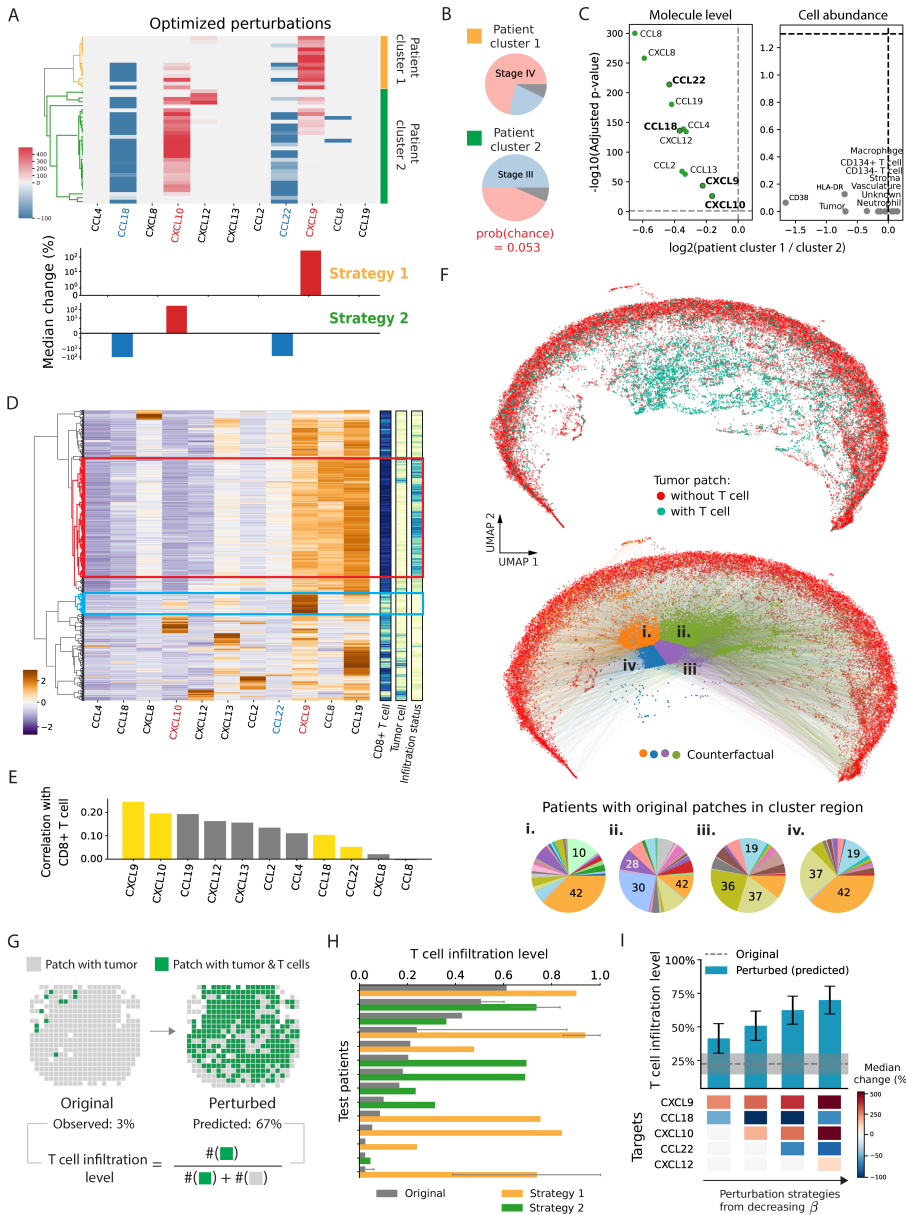


Fig. 3: Combinatorial chemokine therapy predicted to drive T cell infiltration in patients with metastatic melanoma (A) Whole-tumor perturbations optimized across IMC images of patients (row) from the training cohort, with bar graph showing the median relative change in intensity for each molecule.

Fig. 3: (continued) (B) Distribution of cancer stages among patients within two clusters, gray indicates unknown stage, chance probability from hypergeometric distribution. (C) Volcano plot comparing chemokine level and cell type abundance from patient cluster 1 and 2, computed using mean values and Wilcoxon rank sum test. Gray indicates non-statistical significance. (D) Patch-wise chemokine profile (left); 1-D heatmap (right): infiltration status (light/dark = from infiltrated/deserted tumor), tumor cell (light/dark = present/absent), CD8+ T cells (light/dark = present/absent). (E) Patch-wise correlation between chemokine signals and the presence of CD8+ T cells. (F) (Top) UMAP projection of tumor patches (chemokine channels) show a clear separation of masked patches with and without T cells. (Bottom) colored arrows connect UMAP projection of patches without T cells and their corresponding counterfactual (perturbed) patch, where the colors correspond to k-nearest neighbor clusters (i-iv) of the counterfactual patches, highlighting the minimal nature of the perturbations. Pie charts (i-iv) shows the distribution of patients whose original tumor patches are found in the corresponding cluster regions in the UMAP. (G) Cell maps computed from a patient's IMC image, showing the distribution of T cells before and after perturbation. (H) Original vs. perturbed (predicted) mean infiltration level across all patients (test cohort) with 95% confidence interval (only shown for patients with more than 2 samples). Stage IV patients received perturbation strategy 1 (yellow), stage III patients received perturbation strategy 2 (green). (I) Mean infiltration level across all patients (test cohort) for optimized perturbation strategies of varying sparsity, error bar represents 95% CI.

247 Applying our counterfactual optimization procedure using the U-Net classifier
 248 trained on melanoma IMC images, we discovered a combinatorial therapy pre-
 249 dicted to be highly effective in improving T cell infiltration in patients with
 250 melanoma. We restricted the optimization algorithm to only perturb the level
 251 of chemokines, which are a family of secreted proteins that are known for their
 252 ability to stimulate cell migration [31] and have already been harnessed to aug-
 253 ment T-cell therapy [32]. By optimizing over multiple chemokines, Morpheus
 254 opens the door to combinatorial chemokine therapeutics that has the poten-
 255 tially to more effectively enhance T cell infiltration into tumors. **Figure 3A**
 256 shows that patients from the training cohort separate into two clusters based
 257 on hierarchical clustering of perturbations computed for each patient. Tak-
 258 ing median across all patients in cluster 1, the optimized perturbation is to
 259 increase CXCL9 level by 370%, whereas in patient cluster 2, the optimized
 260 perturbation consists of increasing CXCL10 level by 280% while decreasing
 261 CCL18 and CCL22 levels by 100% and 70% respectively (**Figure 3A**). Both
 262 CXCL9 and CXCL10 are well-known for playing a role in the recruitment
 263 of CD8+ T cells to tumors. On the other hand, CCL22 is known to be a
 264 key chemokine for recruiting regulatory T cells [33] and CCL18 is known to
 265 induce an M2-macrophage phenotype [34], so their expression likely promotes

266 an immunosuppressive microenvironment inhibitory to T cell infiltration and
267 function.

268 **Figure 3B** shows that the choice of which of these two strategies was
269 selected for a patient appears to be strongly associated with the patient's cancer
270 stage, with strategy 1 being significantly enriched for patients with stage IV
271 metastatic melanoma and strategy 2 being significantly enriched for patients
272 with stage III cancer, with a probability of 0.053 of such difference being due
273 to chance. Probing deeper into the difference between these two patient clusters,
274 we find that all chemokines have lower mean expression in the tumors of
275 patients in cluster 1 compared to cluster 2, while there are no significant differences
276 between the two groups in terms of the cell type compositions within
277 tumors (**Figure 3C**). Since the levels of CCL22 and CCL18 is 37% and 31%
278 higher in patients from cluster 2 and both chemokines have been implicated in
279 having an inhibitory effect on T-cell infiltration, it is reasonable that the optimization
280 algorithm suggests inhibiting CCL18 and CCL22 only for patients
281 in cluster 2. However, the switch from boosting CXCL9 to CXCL10 is not as
282 straightforward. A possible explanation is that boosting CXCL10 is important
283 when blocking CCL18 and CCL22 in order for the perturbed patches to
284 stay close to the data manifold, leading to more realistic tissue environments.

285 Morpheus selected perturbations that would make the chemokine composition
286 of a TME more similar to T cell rich regions of immune-infiltrated tumors.
287 **Figure 3D** shows that melanoma tissue patches can be clustered into distinct
288 groups based on their chemokine concentration profile. One cluster (highlighted in blue)
289 contains exactly the patches from immune-infiltrated tumors that contain both tumor
290 and T cells, which likely represents a chemokine signature that is suitable for T cell
291 infiltration. Alternately, a second cluster (highlighted in red) which contains patches
292 from immune-desert tumors that have tumor cells but no T cells likely represents an
293 unfavorable chemokine signature. In comparison to the cluster highlighted in red, **Figure 3D**
294 shows the cluster highlighted in blue contains elevated levels of CXCL9, CXCL10 and
295 reduced levels of CCL22 which partially agrees with the perturbation strategy
296 (**Figure 3A**) discovered by Morpheus. Lastly, **Figure 3E** shows that our four
297 selected chemokine targets cannot simply be predicted from correlation of chemokine
298 levels with the presence of CD8+ T cells, as both CCL18 and CCL22 are weakly
299 correlated (< 0.1) with CD8+ T cells even though the optimized perturbations
300 requires inhibiting both chemokines, suggesting the presence of significant nonlinear
301 effects not captured by correlations alone.

302 We can directly observe how Morpheus searches for efficient perturbations
303 by viewing both the original patch and perturbed patches in a dimensionally-reduced
304 space. **Figure 3F** (top) shows a UMAP projection where each point represents the
305 chemokine profile of an IMC patch. T-cell patches (with their CD8+ T cells masked)
306 are well-separated from patches without CD8+ T cells. The colored arrows in the
307 bottom UMAP of **Figure 3F** illustrate the perturbation for each patch as computed
308 by Morpheus, and demonstrate two key
309

310 features of our algorithm. First, optimized perturbations push patches with-
311 out T cells towards the region in UMAP space occupied by T-cell-infiltrated
312 patches. Second, the arrows in Figure 3C are colored to show that optimized
313 perturbations seem efficient in that patches are perturbed just far enough to
314 land in the desired region of space. Specifically, red points that start out on
315 the right edge end up closer to the right after perturbation (region ii and iii),
316 while points that start on the left/bottom edge end up closer to the left/bot-
317 tom (region i and iv), respectively. We make this observation while noting
318 that UMAP, though designed to preserve the topological structure of the data,
319 is not a strictly distance-preserving transformation [35]. Furthermore, the pie
320 charts (i-iv) are colored by the patient of origin to show the region of space
321 where points are being perturbed to are not occupied by tissue samples from
322 a single patient with highly infiltrated tumor. Rather, these regions consist of
323 tissue samples from multiple patients, suggesting that our optimization pro-
324 cedure can synthesize information from different patients when searching for
325 therapeutic strategies.

326 After applying the second perturbation strategy from Figure 3A *in sil-*
327 *ico* to IMC images of a tumor, Figure 3G shows that T cell infiltration level
328 (defined as the proportion of tumor patches with T cells) is predicted to
329 increase by 20 fold. We applied our two perturbation strategies on patients
330 in our test cohort *in silico* after stratifying by cancer stage, using strategy 1
331 on patients with stage IV melanoma and strategy 2 on patients with stage
332 III melanoma. Figure 3H shows that this predicted improvement holds across
333 nearly all 14 patients from the test group, boosting T cell infiltration level
334 from an average of 23% across samples to a predicted 63% post perturbation.
335 For the three test patients with multiple tumor sections (patient 64, 57, 89),
336 we see small to moderate variation in predicted improvement across samples.

337 The combinatorial nature of our optimized perturbation strategy is crucial
338 to its predicted effectiveness. We systematically explored the importance of
339 combinatorial perturbation by changing parameter β of Equation (4) which
340 adjusts the sparsity of the strategy, where a more sparse strategy means fewer
341 molecules are perturbed. Figure 3I shows that perturbing multiple targets
342 is predicted to be necessary for driving significant T cell infiltration across
343 multiple patients, with the best perturbation strategy involving two targets
344 predicted to generate only 60% of the infiltration level achieved by the best
345 perturbation strategy involving four targets. In conclusion, within the scope
346 of the chemokine targets considered, combinatorial perturbation of the TME
347 appears necessary for improving T cell infiltration in metastatic melanoma.

348 Applying Morpheus to CRC with liver metastases
 349 samples

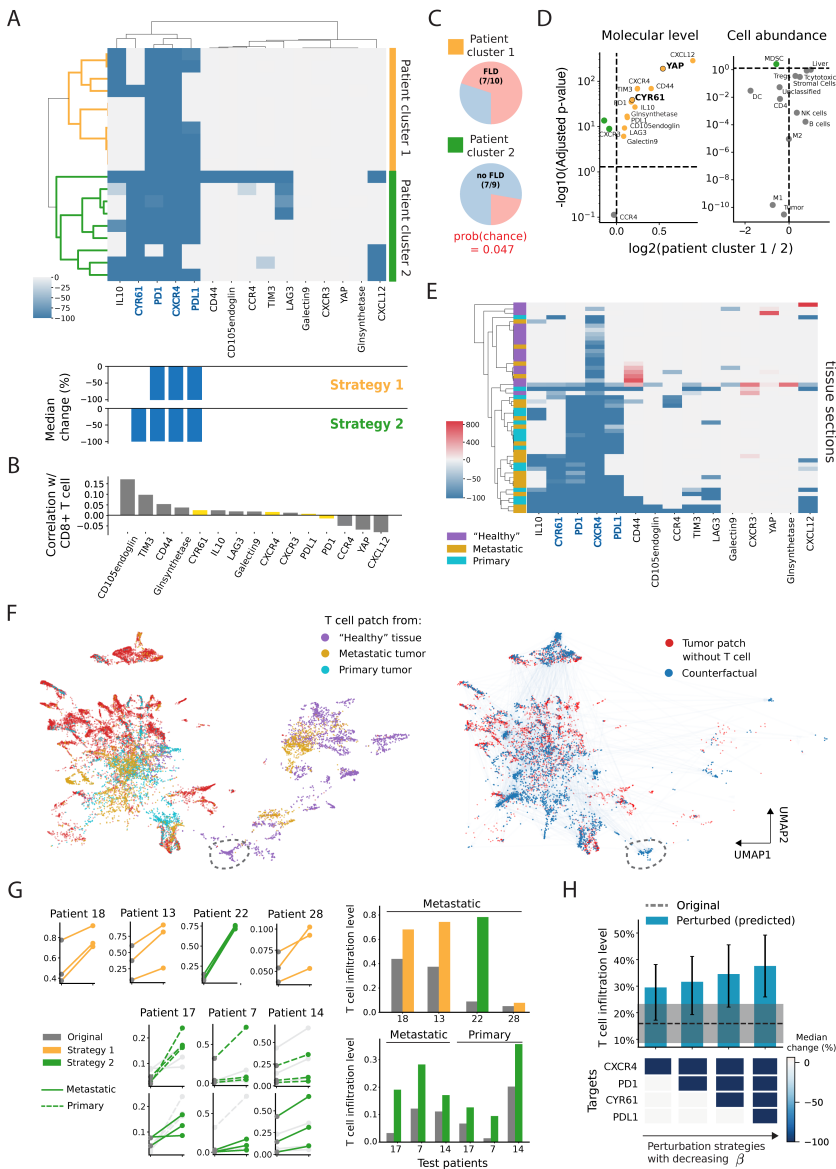


Fig. 4: Blocking subsets of PD-L1, CXCR4, PD-1, and CYR61 predicted to drive T cell infiltration in CRC cohort. (A) Optimized tumor perturbations aggregated to the patient (row) level (train cohort). Bar graph shows the median relative change in intensity for each molecule across all patients within their cluster.

Fig. 4: (continued) (B) Patch-wise correlation between the levels of different molecules and the presence of CD8+ T cells. (C) Pie charts show proportion of patients in each cluster that have fatty liver disease (FLD), chance probability from hypergeometric distribution. (D) Volcano plot comparing molecule levels and cell type abundance between the two patient cluster using tumor tissues, computed using mean values and Wilcoxon rank sum test with Bonferroni correction. (E) Optimized perturbations aggregated to the level of tissue samples (row). (F) UMAP projection of IMC patches, left UMAP shows T cell patches colored by the tissue samples they are taken from. right UMAP shows counterfactual (perturbed) instances optimized for tumor patches without T cells (red). (G) Line plots shows T-cell infiltration level for each tissue section from the test cohort, before and after perturbation. Bar plots show predicted mean T-cell infiltration level for each test patient. (H) Mean infiltration level across all test patients using perturbation strategies of varying sparsity, obtained by varying β in Equation (4), error bar represents 95% CI.

350 Applying Morpheus to IMC images from the CRC cohort, we discovered two
351 patient-dependent therapies predicted to be highly effective in improving T
352 cell infiltration. [Figure 4A](#) shows the optimal perturbations computed for every
353 patient from the training cohort, aggregated over all tumor samples for each
354 patient. Our method consistently discovered two distinct patient-dependent
355 strategies for improving T cell infiltration, as revealed by hierarchical clustering
356 of all patient-level perturbations ([Figure 4A](#)). Taking median over patients in
357 the first cluster, the optimized strategy involves completely inhibiting PD-1,
358 PD-L1, and CXCR4. While for the second group of patients, the optimized
359 strategy involves completely inhibiting CYR61, PD-1, PD-L1, and CXCR4
360 ([Figure 4A](#)). Interestingly, all four of the perturbation targets correlated poorly
361 with the presence of CD8+ T cells compared to the other proteins that were
362 not selected as perturbation targets ([Figure 4B](#)), suggesting the presence of
363 significant spatial and nonlinear effects not captured by correlations alone.

364 All perturbation targets identified by our optimization procedure have been
365 found to play crucial roles in suppressing T cell function in the TME, and treat-
366 ing patients with inhibitors against subsets of the selected targets have already
367 improved T cell infiltration in human CRC liver metastases. Regulatory T
368 cells (Tregs) are recruited into tumor through CXCL12/CXCR4 interaction
369 [36], and the PD-1/PD-L1 pathway inhibits CD8+ T cell activity and infil-
370 tration in tumors. In addition, CYR61 is a chemoattractant and was recently
371 shown to drive M2 TAM infiltration in patients with CRC liver metastases
372 [27]. Inhibition of both PD-1 and CXCR4, which were consistently selected
373 by our algorithm as targets, have already been shown to increase CD8+ T
374 cell infiltration in both patients with CRC and mouse models [37–39]. Finally,
375 [Figure 4A](#) shows that the fifth most common proposed perturbation involves
376 inhibiting IL-10. Indeed, blockade of IL-10 was recently shown to increase the

377 frequency of non-exhausted CD8+ T cell infiltration in slice cultures of human
378 CRC liver metastases [40].

379 The emergence of the two distinct perturbation strategies may be explained
380 by variation in liver fat build-up among patients. Patient cluster 1 is made up
381 of significantly more patients with fatty liver disease (70% FLD) compared to
382 patient cluster 2 (22%), where the probability of this due purely to chance is
383 0.047 (Figure 4C). Furthermore, Figure 4D shows that both YAP and CYR61
384 levels are significantly higher in tumors from patient cluster 1, by 50% and
385 15% respectively. Indeed, CYR61 is known to be associated with non-alcoholic
386 fatty liver disease [27] and YAP is a transcription coregulator that induces
387 CYR61 expression [41]. However despite patients in cluster 1 having higher
388 levels of CYR61, it is only for patients in cluster 2 where the optimal strategy
389 involves blocking CYR61. We postulate that this seemingly paradoxical find-
390 ing may arise because removing CYR61 from patients in cluster 1 represents
391 a more pronounced perturbation, given their inherently higher concentration.
392 A perturbation of this magnitude would likely shift the tumor profile signifi-
393 cantly away from the data manifold, where the classifier’s prediction about the
394 perturbation’s effect becomes less reliable, hence such a perturbation would
395 be heavily penalized during optimization due to the L_{proto} term.

396 Using only raw image patches, Morpheus discovers tissue-dependent per-
397 turbation strategies (Figure 4E). As depicted in Figure 4E, by aggregating
398 perturbations at the individual tissue level, we observe that the optimized
399 perturbation for “healthy” liver sections is straightforward, necessitating only
400 the inhibition of CXCR4. Recall “healthy” sections are samples obtained away
401 from sites of metastasis. In contrast, promoting T cell infiltration into primary
402 colon tumors is anticipated to involve targeting a minimum of three signals.
403 Our method finds that liver metastases appears to fall between these two tissue
404 types. The optimized perturbation strategy for some liver metastases samples
405 is to block CXCR4, while requiring the inhibition of the same set of signals as
406 primary tumors for others. Furthermore, direct comparison between pertur-
407 bations optimized for metastatic tumor and primary tumor samples does not
408 reveal a significant difference in strategy (Figure S2). We can partly under-
409 stand the discrepancy between tissues by plotting a UMAP projection of all
410 T cell patches from the three tissue types (Figure 4F, left). The clear separa-
411 tion between T cell patches from “healthy” tissue and those from primary
412 tumors underscores that the signaling compositions driving T cell infiltration
413 likely differ substantially between the two tissue types. This distinction is
414 likely what prompted our method to identify markedly different perturbation
415 strategies. Furthermore, some patches from metastatic tumors co-localize with
416 “healthy” tissue patches in UMAP space, while other patches co-localizes with
417 primary tumor patches. This observation again aligns with our previous result,
418 where optimized perturbations for metastases samples can bear similarities to
419 strategies for either “healthy” tissue or primary tumor (Figure 4E).

420 Despite the CRC data set comprising a complex blend of healthy, tumor,
421 and hybrid metastatic samples, Morpheus targets the most pertinent tissue

422 type when optimizing perturbations. During both the model training and counter-
423 factual optimization phases, we did not make specific efforts to segregate
424 the three tissue types. Furthermore, we did not provide tissue type labels or
425 any metadata. Despite these nuances, Figure 4F shows that the counterfactual
426 instances for tumor patches (dark blue) from primary and metastases sam-
427 ples are mostly perturbed to be near T cell patches from primary (cyan) and
428 metastatic tumor (gold), instead of being perturbed to be similar to T cell
429 patches from “healthy” tumors (purple). This result is partly a consequence of
430 our prototypical constraint which encourages patches to be perturbed towards
431 the closest T-cell patch. For a patch from a metastatic tumor without T cells,
432 the closest (most similar) T cell patch is likely also from a metastatic tumor
433 than from a “healthy” tissue. However, there are occasional exceptions where
434 T cell patches from “healthy” tissues can influence the optimization of tumor
435 tissues, as outlined by the dashed ellipse in Figure 4F, especially if they share
436 similar features as tumor regions.

437 The two therapeutic strategies we discovered generalize to patients in our
438 test cohort (Figure 4G,H). Given that we have two therapeutic strategies, one
439 enriched for patients with FLD and another for patients without FLD, we
440 apply different perturbation strategies *in silico* across all test patients depend-
441 ing on their FLD status. Aggregated to the patient level, Figure 4G shows that
442 CD8+ T cell infiltration level is predicted to increase for nearly all patients,
443 with the exception of patient 28. Furthermore, aggregating to the entire test
444 cohort, Figure 4H shows a statistically significant boost to mean infiltration
445 level from 15% to a predicted 35% post perturbation. However, when compar-
446 ing individual tissue samples, Figure 4G reveals significant variation in the
447 predicted response to perturbation among samples from the same patient and
448 tissue types. In patient 7, one primary tumor sample is predicted to see a
449 nearly three-fold increase in T cell infiltration after perturbation, yet almost
450 no change is expected for patient 7’s other two primary and three metastatic
451 samples. Similar patterns are observed in patients 14 and 17. This marked
452 variability in response among a significant portion of test patients underscores
453 the challenges posed by intra-tumor and inter-patient heterogeneity in devising
454 therapies for CRC with liver metastases. This result further implies that, for
455 studying CRC with liver metastases, collecting numerous tumor sections per
456 patient could be as crucial as establishing a large patient cohort. Lastly, combi-
457 natorial perturbation is again predicted to be necessary to drive significant
458 T-cell infiltration in large patient cohorts. By increasing β in Equation (4), we
459 generated strategies with between one and four total targets, where our four-
460 target perturbation is the only strategy predicted to produce a statistically
461 significant boost to T-cell infiltration (Figure 4H).

462 Discussion

463 Our integrated deep learning framework, Morpheus, combines deep learn-
464 ing with counterfactual optimization to directly predict therapeutic strategies

465 from spatial omics data. One of the major strengths of Morpheus is that it
466 scales efficiently to deal with large diverse sets of patients samples including
467 metachronous tissue from the same patients but different sites, which will be
468 crucial as more spatial transcriptomics and proteomics data sets are quickly
469 becoming available [42]. Larger data sets could allow us to train more com-
470 plex models such as vision transformers, capturing long range interactions in
471 tissues to improve prediction of T-cell localization. Furthermore, a large set of
472 diverse patient samples will more accurately capture the extent of tumor het-
473 erogeneity, enabling Morpheus to discover therapeutic strategies for different
474 sub-classes of patients.

475 For future work, we would like to apply Morpheus to spatial transcrip-
476 tomics data sets with hundreds to thousands of molecular channels. Although
477 spatial transcriptomics can profile significantly more molecules compared to
478 spatial proteomic techniques [15, 16], the number of spatial transcriptomic
479 profiles of human tumors is currently limited due to the cost, with most pub-
480 lic data sets containing single tissue sections from 1-5 patients which is far
481 too small to apply Morpheus. However, spatial transcriptomics is likely to be
482 more standardized compared to proteomics, which use customized panels. As
483 commercial platforms for spatial transcriptomics start to come online [43], we
484 will likely be seeing large scale spatial transcriptomics data sets in the near
485 future, with ~ 70 -90% of the same probes shared between experiments.

486 A technical extension of Morpheus involves incorporating prior knowl-
487 edge of gene-gene interactions to model the causal relations between genes.
488 Molecular features in tissue profiles can exhibit strong dependencies, there-
489 fore, changing the level of one molecule can affect the expression of others. For
490 example, increased levels of interferon-gamma (IFN- γ) in the tumor microen-
491 vironment, can upregulate the expression of PD-L1 on tumor cells [44]. In
492 order to be more realistic and actionable, a counterfactual should maintain
493 these known causal relations. We can apply a regularizer to penalize counter-
494 factuais that are less feasible according to established gene interactions from
495 knowledge graphs, such as Gene Ontology [45].

496 Other extensions of Morpheus includes predicting cell-type specific pertur-
497 bations, which can be done by directly restricting the perturbation to only
498 alter signals within specific cell types. Additionally, although we applied Mor-
499 pheus to the specific problem of driving T cells to infiltrate solid tumors, we
500 can generalize our framework to predict candidate therapeutics that alter the
501 localization of other cell types. For example, Morpheus can train a classifier
502 model to predict localization of TAMs and compute perturbations predicted
503 to reduce their abundance in the TME.

504 In this work, we focused on identifying generalized therapies by pooling
505 predictions across multiple patient samples, but we can also apply Morpheus
506 to find personalized therapy for treating individual patients. The variation in
507 the optimized perturbations we observe among patients in both melanoma and
508 liver data sets suggest personalize treatments could be significantly more effec-
509 tive compared to generalized therapies (Figure 3A, Figure 4A). Furthermore,

510 **Figure 4G** shows that a therapeutic strategy could have highly variable effect
511 even across different tissue samples from the same patient. This variability sug-
512 gests that to generate therapy for an individual patient, it may be necessary to
513 acquire significant quantities of biopsy data. We can then apply our optimiza-
514 tion procedure to a random subset of the samples, and then test the resulting
515 perturbation strategy on the remaining samples to see how well the strategy
516 is predicted to perform across an entire tumor or other primary/secondary
517 tumors.

518 Incorporating Morpheus in a closed loop with experimental data collection
519 is another promising direction for future work. Data can be collected from
520 patients or animal models with perturbed/engineered signaling context, and
521 this data can be easily fed back into the classifier model to refine the model's
522 prediction. The perturbation could be based on what the model predicts to be
523 effective interventions, as is the case with Morpheus. We can also study tissue
524 samples on which the model tends to make the most mistake and train the
525 model specifically using samples from similar sources, such as similar patient
526 strata or disease state.

527 **Methods**

528 **IMC data sets**

529 All data sets used in this paper are publicly available. Metastatic melanoma
530 data set from Hoch et al. [26] contains 159 images or cores taken from 69
531 patients, collected from sites including skin and lymph-node. CRC liver metas-
532 tases data set from Wang et al. [27] contains 209 images or cores taken from
533 30 patients. Breast tumor data set from Danenberg et al. [28] contains 693
534 images or cores taken from 693 patients. The RNA and protein panels used
535 for each of the three data sets are listed in [Table 1](#).

536 **Data split**

537 For all three IMC data sets, we followed the same data splitting scheme to
538 divide patients into three different groups (training, validation, testing) while
539 ensuring similar class balance across the groups, which in our case means that
540 the proportion of image patches with and without T cells are roughly equal
541 across the three groups for each data set. Specifically, each image within a data
542 set was divided into $48\ \mu\text{m} \times 48\ \mu\text{m}$ patches and the number of patches with
543 and without CD8+ T cells was computed for each image. Furthermore, each
544 patch was downsampled from 48×48 pixels to 16×16 pixel dimension where
545 each pixel now represents a $3\ \mu\text{m} \times 3\ \mu\text{m}$ region. We applied spectral analysis
546 to study the effect of using different patch size to predict T cell infiltration
547 and found that our selected patch size remains highly informative of T cell
548 presence ([Figure S1](#)). Next, the patients are shuffled between the three groups
549 until three criteria are met: 1) the number of patients across the three groups
550 follow a 65/15/20 ratio, 2) the difference in class proportion between any two

Metastatic melanoma		CRC with liver metastases		Breast tumor	
Vimentin	DapB	CD45	Glnsynthetase	Histone H3	SMA
CD163	CCL4	CD163	NKG2D	CK5	CD38
B2M	CCL18	CCR4	PD-L1	HLA-DR	CK8-18
CD134	CXCL8	FAP	CD11c	CD15	FSP1
CD68	CXCL10	LAG3	HepPar1	CD163	ICOS
GLUT1	CXCL12	FOXP3	αSMA	OX40	CD68
CD3	CXCL13	CD4	CD105	HER2 (3B5)	CD3
LAG3	CCL2	CD68	VISTA	Podoplanin	CD11c
PD-1	CCL22	CD20	CD8α	PD-1	GITR
HistoneH3	CXCL9	TIM3	CXCR4	CD16	c-Caspase3
CCR2	CCL19	PD-1	iNOS	CD45RA	B2M
PD-L1	CCL8	CD31	CYR61	CD45RO	FOXP3
CD8	SMA	CDX2	CAIX	CD20	ER
SOX10	CD31	CD3	CD44	CD8	CD57
Mart1	pRB	CD15	CD11b	Ki-67	PDGFR β
cleavedPARP	MPO	HLA-DR	IL10	Caveolin-1	CD4
CD15	CK5	CXCL12	HLA-ABC	CD31-vWF	CXCL12
CD38	HLA-DR	GranzymeB	Ki67	HLA-ABC	panCK
S100	Cadherin11	HistoneH3	CXCR3	HER2 (D8F12)	
FAP		Galectin9	YAP		
		CD14	CK19		

Table 1: Protein and RNA panels imaged for each of the IMC data sets, with RNA targets bolded

551 of the three groups is less than 2%, and 3) the training set contains at least
552 65% of total patches. The actual data splits used in the paper are described
553 in [Table 2](#).

Data set	Group	Patient count	Patch count	Proportion of patches with CD8+ T cells
Metastatic melanoma	Training	102	23741	29.6%
	Validation	28	6045	30.3%
	Testing	29	5950	30.4%
CRC with liver metastases	Training	19	44449	15.9%
	Validation	4	6957	14.4%
	Testing	7	14907	15.9%
Breast cancer	Training	485	41104	23.7%
	Validation	113	9015	23.4%
	Testing	151	12987	23.8%

Table 2: Data split for Melanoma, CRC cohort, and breast tumor IMC data set

Single-cell phenotyping

For each data set, we used the cell type classification (tumor and CD8+ T cells) from the original paper. For the melanoma data set, cell phenotyping was performed using the Shiny application of the R package cytomapper [46], which allows labeling of cell populations using multiple gates. CD8+ T cells were defined using CD3 and CD8, tumor cells are positive for any or multiple of SOX9, SOX10, MITF, Mart1, S100A1, and p75. For the CRC and breast cancer data set, cell type labeling was performed using PhenoGraph [47].

Classifier training

In this work, we trained three classes of models to perform our T cell prediction task. All models presented in this paper were trained with early stopping based on the validation Matthews Correlation Coefficient (MCC) for 10-20 epochs. All models were trained on an NVIDIA GeForce RTX 3090 Ti GPU using PyTorch version 1.13.1 [48]. More details about hyperparameters and implementations can be found in our Github repository.

T cell masking strategy

The purpose of model training is for the model to learn molecular features of the CD8+ T cell's environment that is indicative of its presence, so it is important for us to remove features of the image that are predictive of CD8+ T cell presence but are not part of the cell's environment. We devised a non-trivial cell masking strategy in order to remove T-cell expression patterns without introducing new features that are highly predictive of T cell presence but are not biologically relevant. A simple masking strategy of zeroing out all pixels belonging to CD8+ T cells will introduce contiguous regions of zeros to image patches with T cells, which is an artificial feature that is nonetheless highly predictive of T-cell presence and hence will likely be the main feature learned by a model during training. To circumvent this issue, we first apply a cell "pixelation" step to the original IMC image where we reduce each cell to a single pixel positioned at the cell's centroid. The value of this pixel is the sum of all pixels originally associated with the cell, representing the total signal from each channel within the cell. We then mask this "pixelated" image by zeroing all pixels representing CD8+ T cells. Since there are usually at most two T cell pixels in an image patch, zeroing them in a 16×16 pixel image where most ($> 90\%$) of the pixels are already zeroes is not likely to introduce a significant signal that is predictive T cell presence. We show that our strategy is effective at masking T cells without introducing additional features through a series of image augmentation experiments ([Supplemental Note 1 Assessment of T-cell masking strategy](#)).

Logistic regression models

We trained a single-layer neural network on the average intensity values from each molecular channel to obtain a logistic regression classifier, predicting the

595 probability of CD8+ T cell presence in the image patch. This model represents
 596 a linear model where only the average intensity of each molecule is used for
 597 prediction instead of their spatial distribution within a patch.

598 MLP models

599 Similar to a logistic regression model, the Multilayer Perceptron (MLP) also
 600 uses averaged intensity as input features for prediction but is capable of learn-
 601 ing nonlinear interactions between features. The MLP model consists of two
 602 hidden layers (30 and 10 nodes) with ReLU activation.

603 U-Net models

604 To train networks that can make full use of the spatial information, we used
 605 a fully convolutional neural network with the U-Net architecture. The U-Net
 606 architecture consists of a contracting path and an expansive path, which gives
 607 it a U-shaped structure [29]. The contracting path consists of four repeated
 608 blocks, each containing a convolutional layer followed by a Rectified Lin-
 609 ear Unit (ReLU) activation and a max pooling layer. The expansive path
 610 mirrors the contracting path, where each block contains a transposed convolu-
 611 tional layer. Skip connections concatenates the up-sampled features with the
 612 corresponding feature maps from the contracting path to include local infor-
 613 mation. The output of the expansive path is then fed to a fully-connected layer
 614 with softmax activation to produce a predicted probability. The model was
 615 trained from scratch using image augmentation to prevent over-fitting, includ-
 616 ing random horizontal/vertical flips and rotations, in addition to standard
 617 channel-wise normalization. We train our U-net classifiers using stochastic gra-
 618 dient descent with momentum and a learning rate of 10^{-2} on mini-batches of
 619 size 128.

620 Counterfactual optimization

621 Given an IMC patch $x^{(i)}$ without T cells, and a classifier f , our goal is to find
 622 a perturbation $\delta^{(i)}$ for the patch such that f classifies the perturbed patch as
 623 having T cells. For CNN models, $\delta^{(i)} \in \mathbb{R}^{u \times l \times d}$ is a 3D tensor that describes
 624 changes made for every channel, at each pixel of the patch.

625 Given a CNN classifier f and a IMC patch $x^{(i)}$ such that $f(x_0^{(i)}) =$
 626 $\mathbb{P}(\text{T cells present}) < p$, where $p > 0$ is the classification threshold below which
 627 the classifier predicts no T-cell, we aim to obtain a perturbation $\delta^{(i)}$ such that
 628 $f(x_0^{(i)} + \delta^{(i)}) > p$, by solving the following optimization problem adopted from
 629 [30],

$$\delta^{(i)} = \min_{\delta} L_{\text{pred}}(x_0^{(i)}, \delta) + L_{\text{dist}}(\delta) + L_{\text{proto}}(x_0^{(i)}, \delta), \quad (7)$$

such that

$$L_{\text{pred}}(x_0^{(i)}, \delta) = c \max(-f(x_0^{(i)} + \delta), -p), \quad (8)$$

$$L_{\text{dist}}(\delta) = \beta \|\delta\|_1 + \|\delta\|_2^2, \quad (9)$$

$$L_{\text{proto}}(x_0^{(i)}, \delta) = \theta \|x_0^{(i)} + \delta - \text{proto}^{(i)}\|_2^2, \quad (10)$$

$$\delta^{(i)} = \gamma^{(i)} \odot_3 x_0^{(i)} \quad (11)$$

630 where $\text{proto}^{(i)}$ is an instance of the training set classified as having T cells,
 631 defined by first building a k-d tree of training instances classified as having T
 632 cells and setting the k -nearest item in the tree (in terms of euclidean distance
 633 to $x_0^{(i)}$) as proto. We use $k = 1$ for all counterfactual optimization. For all other
 634 parameters, we list their values in [Table 3](#). During optimization, the weight c
 635 of the loss term L_{pred} is updated for n iterations, starting at c_0 . If we identify
 636 a valid counterfactual for the present value of c , we will then decrease c in
 637 the subsequent optimization cycle to increase the weight of the additional loss
 638 components, thereby enhancing the overall solution. If, however, we do not
 639 identify a counterfactual, c is increased to put more emphasis on increasing
 640 the predicted probability of the counterfactual. The parameter s_{max} sets the
 641 maximum number of optimization steps for each value of c .

Parameters	Melanoma	CRC
β	2	80
θ	60	40
p	0.5	0.43
c_0	1000	1000
n	5	5
s_{max}	1000	1000

Table 3: Parameter values used for counterfactual optimization

642 Code Availability

643 Code for model training, perturbation optimization and analysis are publicly
 644 available at <https://github.com/neonine2/morpheus>. Our optimization code
 645 was implemented in Python and was built upon the open source Python library
 646 Alibi [49].

647 Data Availability

648 All data sets used in this study are published and publicly available.

649 Acknowledgements

We would like to thank Inna Strazhnik for her support with figure illustrations. We would like to thank Akil Merchant, Alma Andersson, Aviv Regev, Long Cai, Barbara Wold, Michal Polonsky, Jonathan Fox, Yujing Yang, Abdullah Farooq and all members of the Thomson lab for insightful discussion that

significantly improved this work. We gratefully acknowledge the support of the National Institutes of Health's Information Technology for Cancer Research (ITCR) program and the Merkin Institute for Translational Research.

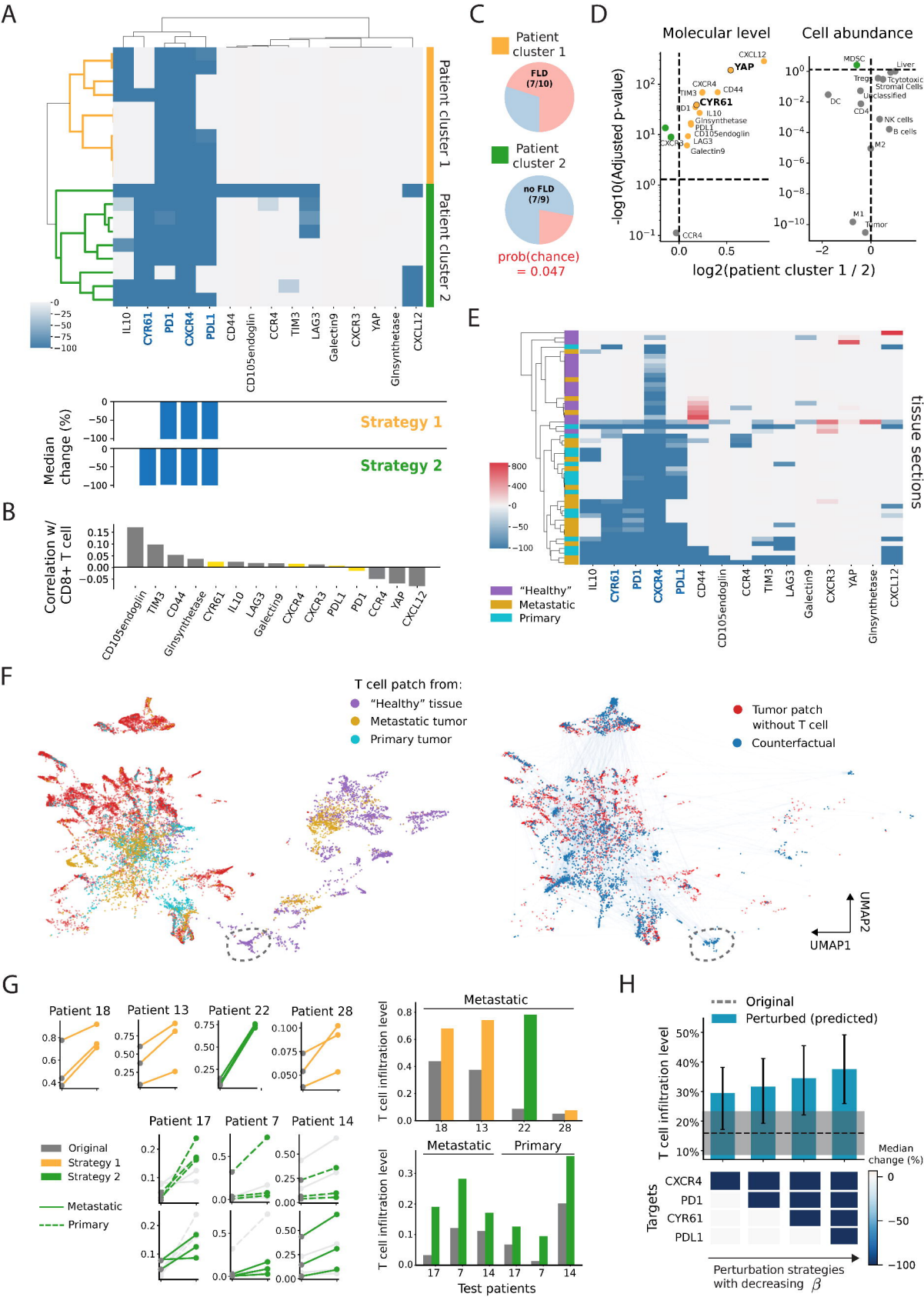
References

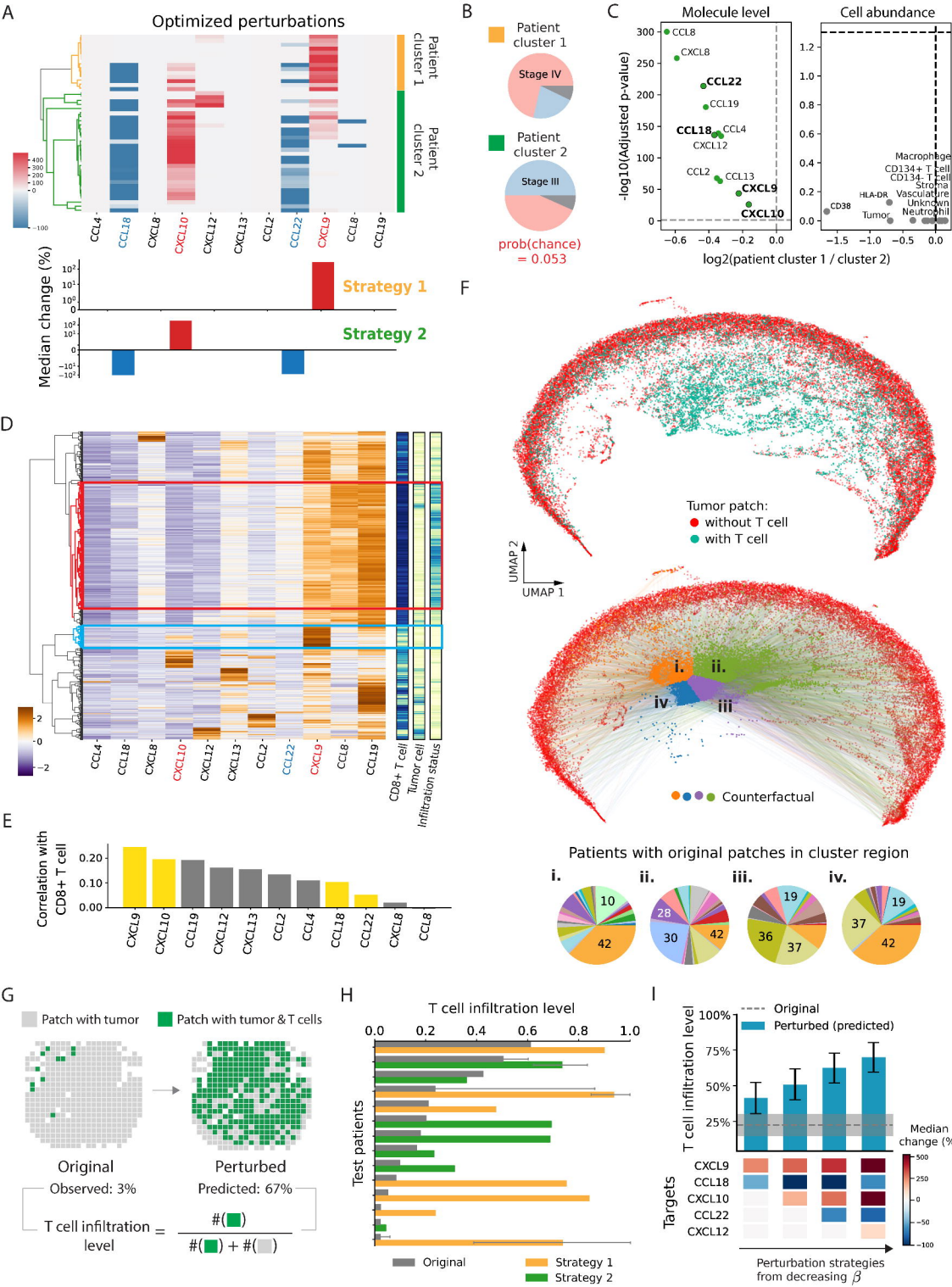
- [1] Fridman, W. H., Zitvogel, L., Sautès-Fridman, C. & Kroemer, G. The immune contexture in cancer prognosis and treatment. *Nature reviews Clinical oncology* **14** (12), 717–734 (2017) .
- [2] Binnewies, M. *et al.* Understanding the tumor immune microenvironment (time) for effective therapy. *Nature medicine* **24** (5), 541–550 (2018) .
- [3] Bruni, D., Angell, H. K. & Galon, J. The immune contexture and immunoscore in cancer prognosis and therapeutic efficacy. *Nature Reviews Cancer* **20** (11), 662–680 (2020) .
- [4] Hegde, P. S. & Chen, D. S. Top 10 challenges in cancer immunotherapy. *Immunity* **52** (1), 17–35 (2020) .
- [5] Choe, J. H., Williams, J. Z. & Lim, W. A. Engineering t cells to treat cancer: the convergence of immuno-oncology and synthetic biology. *Annual Review of Cancer Biology* **4**, 121–139 (2020) .
- [6] Pitt, J. *et al.* Targeting the tumor microenvironment: removing obstruction to anticancer immune responses and immunotherapy. *Annals of Oncology* **27** (8), 1482–1492 (2016) .
- [7] Haslam, A. & Prasad, V. Estimation of the percentage of us patients with cancer who are eligible for and respond to checkpoint inhibitor immunotherapy drugs. *JAMA network open* **2** (5), e192535–e192535 (2019) .
- [8] Lee, J. S. & Ruppin, E. Multiomics Prediction of Response Rates to Therapies to Inhibit Programmed Cell Death 1 and Programmed Cell Death 1 Ligand 1. *JAMA oncology* **5** (11), 1614–1618 (2019) .
- [9] Pittet, M. J., Michielin, O. & Migliorini, D. Clinical relevance of tumour-associated macrophages. *Nature reviews Clinical oncology* **19** (6), 402–421 (2022) .
- [10] Bonaventura, P. *et al.* Cold tumors: a therapeutic challenge for immunotherapy. *Frontiers in immunology* **10**, 168 (2019) .
- [11] Savas, P. *et al.* Clinical relevance of host immunity in breast cancer: from TILs to the clinic. *Nature reviews Clinical oncology* **13** (4), 228–241 (2016) .
- [12] Tsauro, I., Brandt, M. P., Juengel, E., Manceau, C. & Ploussard, G. Immunotherapy in prostate cancer: new horizon of hurdles and hopes. *World journal of urology* **39**, 1387–1403 (2021) .

- [13] Moffitt, J. R., Lundberg, E. & Heyn, H. The emerging landscape of spatial profiling technologies. *Nature Reviews Genetics* **23** (12), 741–759 (2022) .
- [14] Lanitis, E., Dangaj, D., Irving, M. & Coukos, G. Mechanisms regulating t-cell infiltration and activity in solid tumors. *Annals of Oncology* **28**, xii18–xii32 (2017) .
- [15] Rodriques, S. G. *et al.* Slide-seq: A scalable technology for measuring genome-wide expression at high spatial resolution. *Science* **363** (6434), 1463–1467 (2019) .
- [16] Eng, C.-H. L. *et al.* Transcriptome-scale super-resolved imaging in tissues by RNA seqFISH+. *Nature* **568** (7751), 235–239 (2019) .
- [17] Giesen, C. *et al.* Highly multiplexed imaging of tumor tissues with sub-cellular resolution by mass cytometry. *Nature methods* **11** (4), 417–422 (2014) .
- [18] Goltsev, Y. *et al.* Deep profiling of mouse splenic architecture with CODEX multiplexed imaging. *Cell* **174** (4), 968–981 (2018) .
- [19] Bhate, S. S., Barlow, G. L., Schürch, C. M. & Nolan, G. P. Tissue schematics map the specialization of immune tissue motifs and their appropriation by tumors. *Cell Systems* **13** (2), 109–130 (2022) .
- [20] Wu, Z. *et al.* Graph deep learning for the characterization of tumour microenvironments from spatial protein profiles in tissue specimens. *Nature Biomedical Engineering* 1–14 (2022) .
- [21] Schürch, C. M. *et al.* Coordinated Cellular Neighborhoods Orchestrate Antitumoral Immunity at the Colorectal Cancer Invasive Front. *Cell* **182** (5), 1341–1359 (2020) .
- [22] Aoki, T. *et al.* The spatially resolved tumor microenvironment predicts treatment outcome in relapsed/refractory Hodgkin lymphoma. *bioRxiv* 2023–05 (2023) .
- [23] Chang, C.-H., Creager, E., Goldenberg, A. & Duvenaud, D. Explaining Image Classifiers by Counterfactual Generation. *International Conference on Learning Representations (ICLR)* (2019) .
- [24] Wachter, S., Mittelstadt, B. & Russell, C. Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harv. JL & Tech.* **31**, 841 (2017) .
- [25] Verma, S. *et al.* Counterfactual explanations and algorithmic recourses for machine learning: A review. *arXiv preprint arXiv:2010.10596* (2020) .

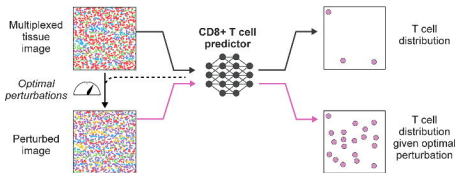
- [26] Hoch, T. *et al.* Multiplexed imaging mass cytometry of the chemokine milieu in melanoma characterizes features of the response to immunotherapy. *Science Immunology* **7** (70) (2022) .
- [27] Wang, Z. *et al.* Extracellular vesicles in fatty liver promote a metastatic tumor microenvironment. *Cell Metabolism* (2023) .
- [28] Danenberg, E. *et al.* Breast tumor microenvironment structures are associated with genomic features and clinical outcome. *Nature genetics* **54** (5), 660–669 (2022) .
- [29] Buda, M., Saha, A. & Mazurowski, M. A. Association of genomic subtypes of lower-grade gliomas with shape features automatically extracted by a deep learning algorithm. *Computers in Biology and Medicine* **109** (2019). <https://doi.org/10.1016/j.combiomed.2019.05.002> .
- [30] Looveren, A. V. & Klaise, J. Interpretable counterfactual explanations guided by prototypes. *Joint European Conference on Machine Learning and Knowledge Discovery in Databases* 650–665 (2021) .
- [31] Hughes, C. E. & Nibbs, R. J. A guide to chemokines and their receptors. *The FEBS journal* **285** (16), 2944–2971 (2018) .
- [32] Foeng, J., Comerford, I. & McColl, S. R. Harnessing the chemokine system to home CAR-T cells into solid tumors. *Cell Reports Medicine* (2022) .
- [33] Kohli, K., Pillarisetty, V. G. & Kim, T. S. Key chemokines direct migration of immune cells in solid tumors. *Cancer gene therapy* **29** (1), 10–21 (2022) .
- [34] Schraufstatter, I. U., Zhao, M., Khaldoyanidi, S. K. & DiScipio, R. G. The chemokine CCL18 causes maturation of cultured monocytes to macrophages in the M2 spectrum. *Immunology* **135** (4), 287–298 (2012) .
- [35] McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *ArXiv e-prints* (2018). <https://arxiv.org/abs/1802.03426> .
- [36] Ghanem, I. *et al.* Insights on the CXCL12-CXCR4 axis in hepatocellular carcinoma carcinogenesis. *American journal of translational research* **6** (4), 340 (2014) .
- [37] Biasci, D. *et al.* CXCR4 inhibition in human pancreatic and colorectal cancers induces an integrated immune response. *Proceedings of the National Academy of Sciences* **117** (46), 28960–28970 (2020) .

- [38] Chen, Y. *et al.* CXCR4 inhibition in tumor microenvironment facilitates anti-programmed death receptor-1 immunotherapy in sorafenib-treated hepatocellular carcinoma in mice. *Hepatology* **61** (5), 1591–1602 (2015) .
- [39] Steele, M. M. *et al.* T cell egress via lymphatic vessels is tuned by antigen encounter and limits tumor control. *Nature Immunology* **24** (4), 664–675 (2023) .
- [40] Sullivan, K. M. *et al.* Blockade of interleukin 10 potentiates antitumour immune function in human colorectal cancer liver metastases. *Gut* **72** (2), 325–337 (2023) .
- [41] Zhang, H., Pasolli, H. A. & Fuchs, E. Yes-associated protein (YAP) transcriptional coactivator functions in balancing growth and differentiation in skin. *Proceedings of the National Academy of Sciences* **108** (6), 2270–2275 (2011) .
- [42] Chen, A. *et al.* Spatiotemporal transcriptomic atlas of mouse organogenesis using dna nanoball-patterned arrays. *Cell* **185** (10), 1777–1792 (2022) .
- [43] Janesick, A. *et al.* High resolution mapping of the breast cancer tumor microenvironment using integrated single cell, spatial and in situ analysis of FFPE tissue. *Biorxiv* 2022–10 (2022) .
- [44] Qian, J. *et al.* The IFN- γ /PD-L1 axis between T cells and tumor microenvironment: hints for glioma anti-PD-1/PD-L1 therapy. *Journal of neuroinflammation* **15** (1), 1–13 (2018) .
- [45] Consortium, G. O. The Gene Ontology (GO) database and informatics resource. *Nucleic acids research* **32**, D258–D261 (2004) .
- [46] Eling, N., Damond, N., Hoch, T. & Bodenmiller, B. cytomapper: an R/Bioconductor package for visualization of highly multiplexed imaging data. *Bioinformatics* **36** (24), 5706–5708 (2020) .
- [47] Levine, J. H. *et al.* Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells that Correlate with Prognosis. *Cell* **162** (1), 184–197 (2015) .
- [48] Paszke, A. *et al.* PyTorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* **32** (2019) .
- [49] Klaise, J., Looveren, A. V., Vacanti, G. & Coca, A. Alibi Explain: Algorithms for Explaining Machine Learning Models. *Journal of Machine Learning Research* **22** (181), 1–7 (2021) .

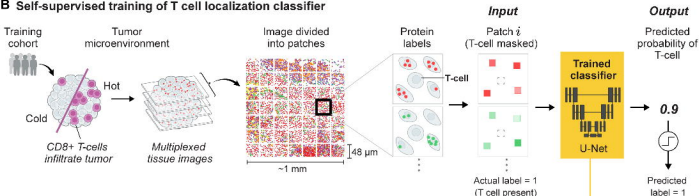




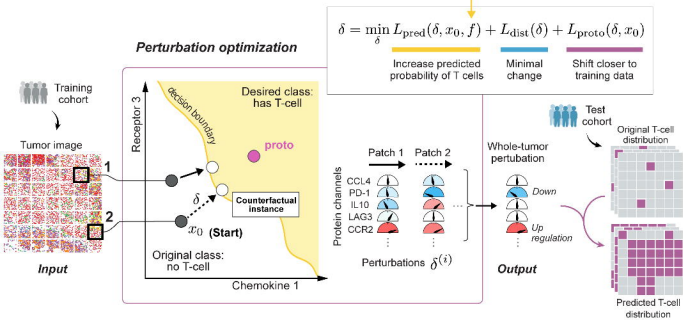
A Overview of Morpheus: a counterfactual optimization framework



B Self-supervised training of T cell localization classifier

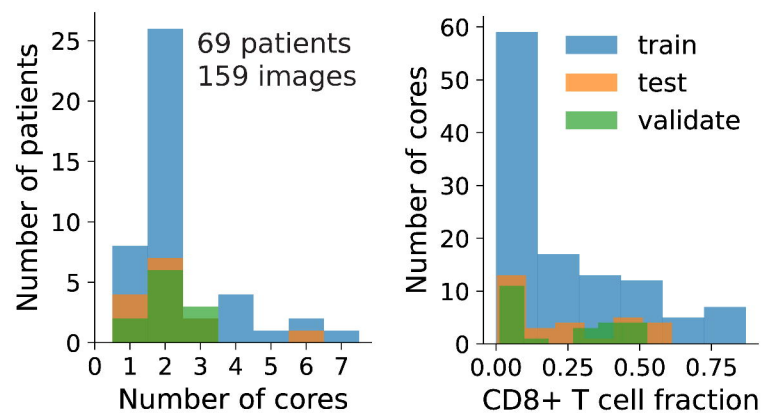


C Counterfactual optimization of tissue perturbation

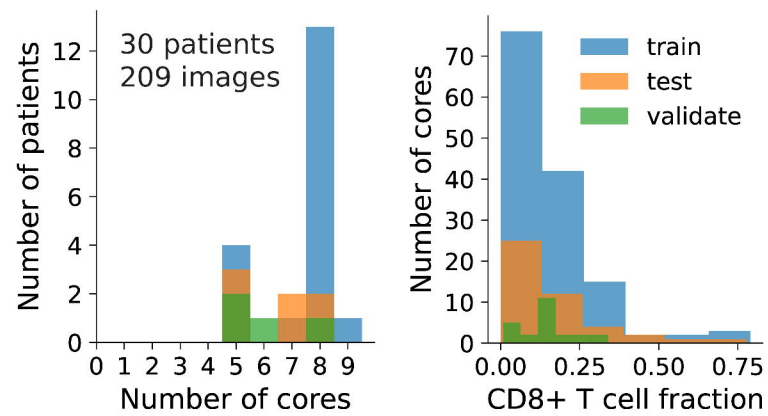


A

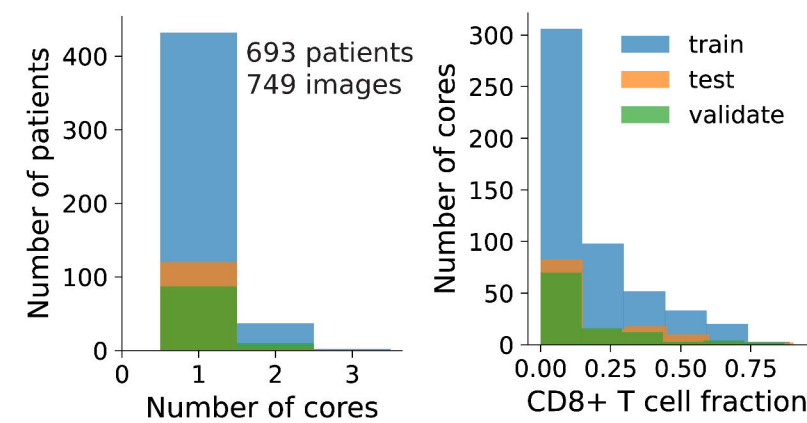
Melanoma



CRC with liver metastases

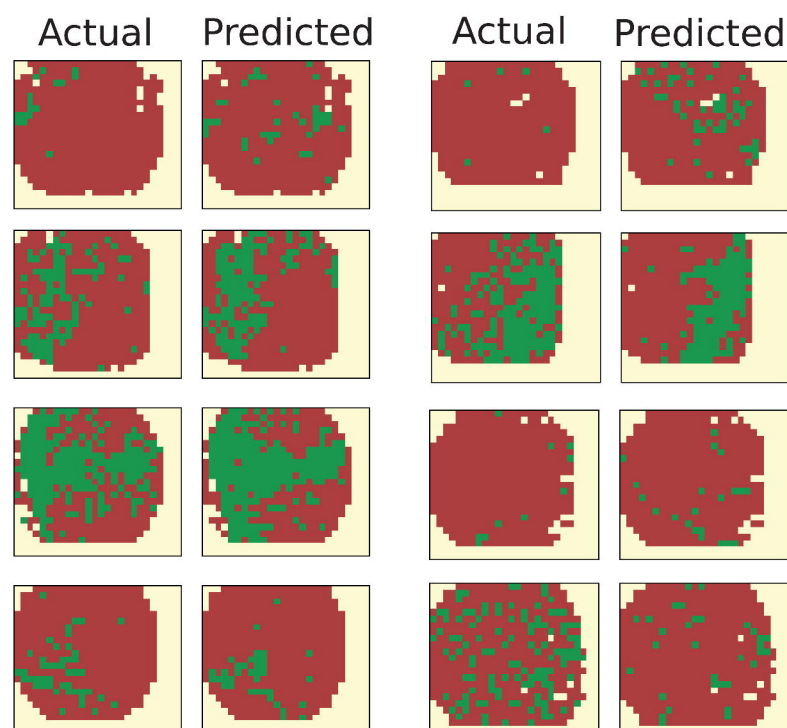


Breast tumor

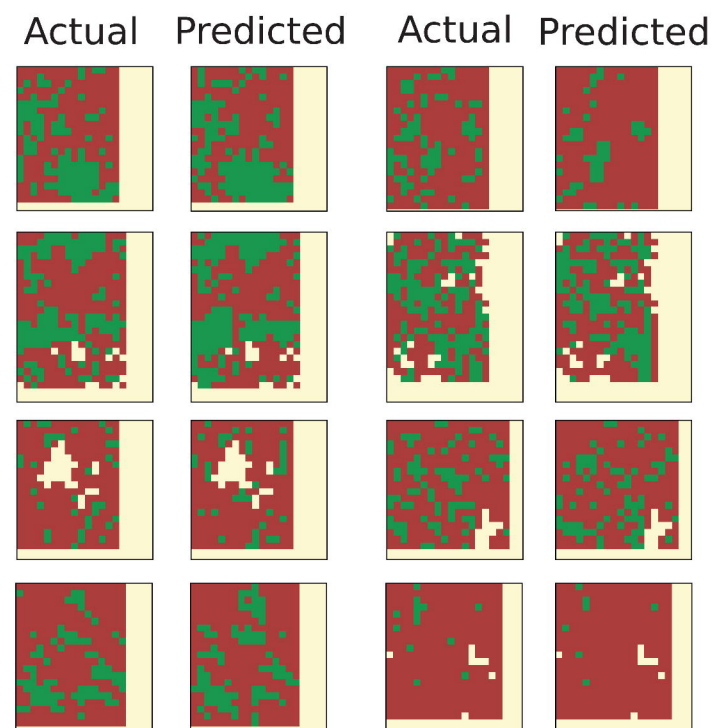


B

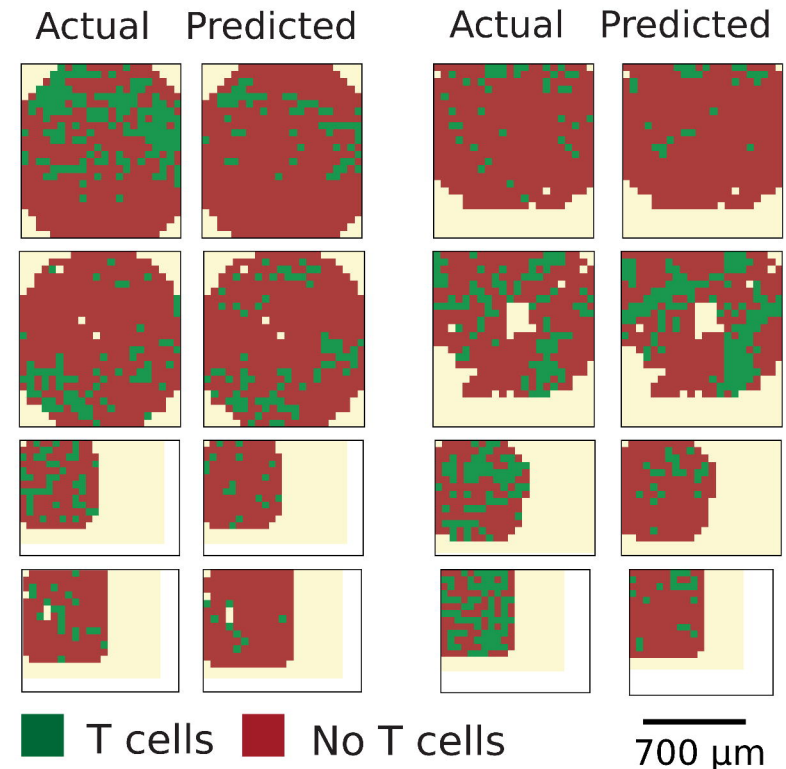
Melanoma



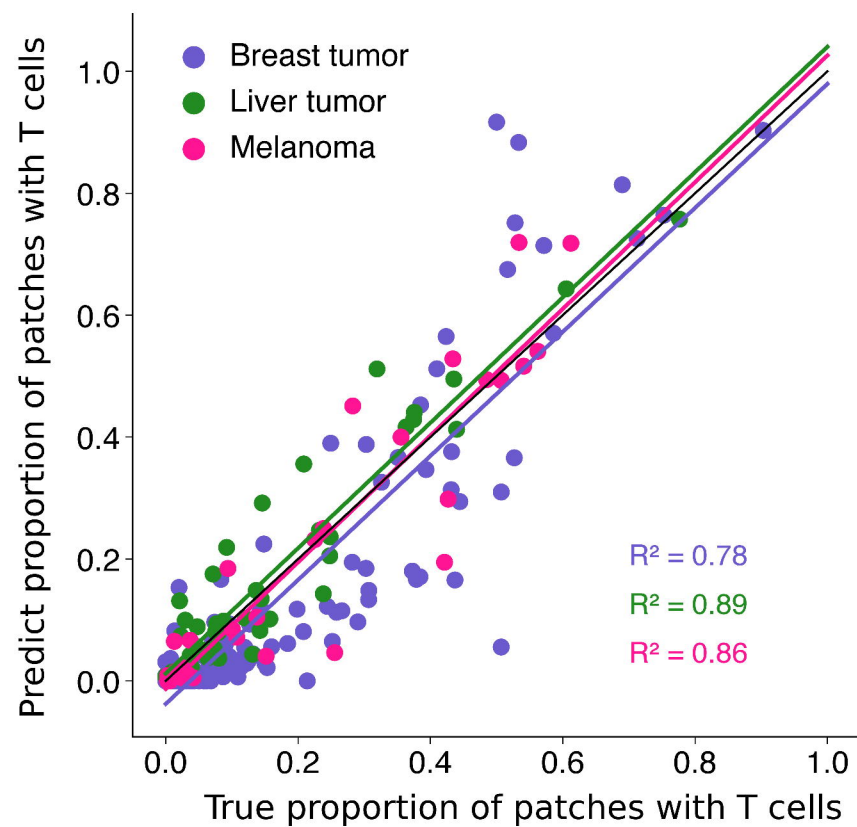
CRC with liver metastases



Breast tumor



C



D

